

CONVEX CATEGORIES AND ENTROPY: NOTES FOR UCONN'S S.I.G.M.A. SEMINAR

ARTHUR J. PARZYG NAT

ABSTRACT. The concept of entropy arises in a plethora of disparate scenarios. To name a few, there are Boltzmann's definition in terms of statistical ensembles of particles, Shannon's information-theoretic notion of entropy, von Neumann's (quantum) entropy, Bekenstein's black hole entropy, Kolmogorov-Sinai's entropy for dynamical systems, and even Perelman's entropy in geometry which was an ingredient in Perelman's proof of the Poincaré conjecture. This talk focuses on Shannon entropy from a categorical perspective and provides an overview of recent work of Baez, Fritz, and Leinster. The main point of their work is that although entropy is not characterized as a convex function on probability spaces, it is characterized as a convex functor on the category of probability spaces and probability-preserving processes. All categorical background will be supplied.

CONTENTS

1. Entropy and information as uncertainty—some early theorems	1
2. Other instances of entropy	4
3. The category of finite probability spaces	6
4. The convex structure on the category of finite probability spaces	9
5. Entropy is an affine functor	13
References	15

1. ENTROPY AND INFORMATION AS UNCERTAINTY—SOME EARLY THEOREMS

Definition 1.1. A finite probability space is a pair (X, p) consisting of a finite set X and a function $p : X \rightarrow \mathbb{R}$ satisfying

$$(1.2) \quad p(x) \geq 0 \quad \forall x \in X$$

and

$$(1.3) \quad \sum_{x \in X} p(x) = 1.$$

The elements of X are called events and p is called a probability distribution on X .

Date: December 2, 2016.

Example 1.4. Let $X = \{H, T\}$ be a two-element set (H stands for “heads” and T stands for “tails”) and let p be the constant function $\frac{1}{2}$. Then (X, p) is interpreted as a fair coin toss. Let $\lambda \in [0, 1]$ and define $p_\lambda : X \rightarrow \mathbb{R}$ to be

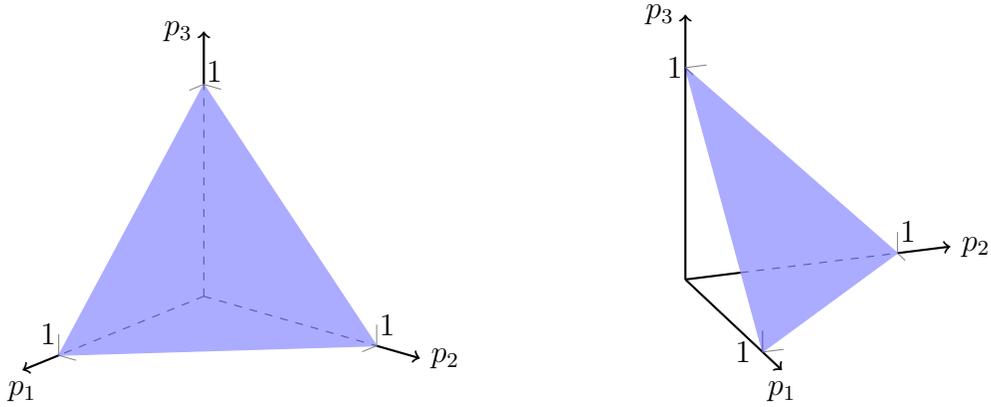
$$(1.5) \quad p_\lambda(H) := \lambda \quad \& \quad p_\lambda(T) := 1 - \lambda.$$

Then (X, p_λ) is interpreted as an unfair coin toss whenever $\lambda \neq \frac{1}{2}$. Both situations are referred to as a coin toss. Let $\mathbf{2}$ denote the fair coin toss.

The set of probability distributions on a finite set X of cardinality $n := |X|$ is in one-to-one correspondence with elements of an $(n - 1)$ -simplex Δ^{n-1} , which is defined to be

$$(1.6) \quad \Delta^{n-1} := \left\{ (p_1, \dots, p_n) \in \mathbb{R}^n \mid \sum_{i=1}^n p_i = 1 \text{ and } p_i \geq 0 \forall i = 1, \dots, n \right\}.$$

For example, the 2-simplex looks like the following subset of \mathbb{R}^3 viewed from two different angles



Note, also, that the permutation group S_n acts on \mathbb{R}^n by permutation of coordinates and this action induces one on Δ^{n-1} for all $n \in \mathbb{N}$.

Definition 1.7. Let Δ denote the disjoint union of all simplices,

$$(1.8) \quad \Delta := \coprod_{n=1}^{\infty} \Delta^{n-1}.$$

Theorem 1.9 (Fadeev’s characterization [13]). *Let $H : \Delta \rightarrow \mathbb{R}$ be a function satisfying the following four axioms.*

(a) *For each $n \in \mathbb{N}$, the diagram*

$$(1.10) \quad \begin{array}{ccc} \Delta^{n-1} & \xrightarrow{\sigma} & \Delta^{n-1} \\ & \searrow H & \swarrow H \\ & & \mathbb{R} \end{array}$$

commutes for all $\sigma \in S_n$.

(b) The function $H : \Delta^1 \rightarrow \mathbb{R}$ is continuous.

(c) $H(\frac{1}{2}, \frac{1}{2}) = 1$.

(d) For each $n \in \mathbb{N}$,

$$(1.11) \quad H(tp_1, (1-t)p_1, p_2, \dots, p_n) = H(p_1, \dots, p_n) + p_1 H(t, 1-t)$$

for all $t \in [0, 1]$.

Then

$$(1.12) \quad H(p_1, \dots, p_n) = - \sum_{k=1}^n p_k \log_2 p_k.$$

Physics 1.13. The interpretation of the function H is the measure of uncertainty or the information content of a probability distribution. Fadeev's axioms can therefore be interpreted in the following way.

- (a) Relabeling equally probable events does not change the uncertainty of the system.
- (b) Infinitesimally perturbing the unfair coin toss results in an infinitesimal perturbation of the uncertainty.
- (c) This is a normalization to specify the elementary unit for measuring uncertainty or information. The amount of uncertainty or information in the fair coin toss is declared to be 1.
- (d) Consider a set $X = \{x_1, \dots, x_n\}$ of n events with probabilities p_1, \dots, p_n , respectively. Duplicate the event x_1 and label the two copies y_1 and z_1 so that they are distinguishable. Then, apply a coin toss of weight $\lambda \in [0, 1]$ to choose one of these two copies. The resulting set consists of $n + 1$ elements $\{y_1, z_1, x_2, \dots, x_n\}$ and the probability distribution $\lambda p_1, (1-\lambda)p_1, p_2, \dots, p_n$, respectively. The measure of uncertainty for this new probability distribution is the uncertainty in the first probability distribution plus the uncertainty of the coin toss multiplied by the probability of the event x_1 of occurring.

One can show from these axioms that the maximal measure of uncertainty on n elements occurs for the probability distribution $p_k = \frac{1}{n}$ for all $k \in \{1, \dots, n\}$.

Definition 1.14. The function H from Fadeev's Theorem is known as the Shannon entropy [14]. If (X, p) is a finite probability space, then

$$(1.15) \quad H(p) := - \sum_{x \in X} p(x) \log_2 (p(x))$$

is the Shannon entropy of (X, p) .

If you thought that last axiom from Fadeev's Theorem was difficult to understand, then you might appreciate more intuitive yet mathematically faithful axiomatizations for the Shannon entropy, i.e. the measure of uncertainty of a probability distribution. Describing one such axiomatization is our goal.

2. OTHER INSTANCES OF ENTROPY

The notion of entropy appears in several different contexts yet shares many similar features. Probably the most ubiquitous property that is common to all forms is concavity (which I will refer to as convexity throughout this discussion). This is because entropy is often interpreted as the lack of information or ignorance. And if you take two systems and combine them in a convex fashion, then you lose the information that distinguishes from which of the subsystems a state comes from [16]. Thus, one expects the ignorance to be larger than the weighted sum of the individual entropies.

Example 2.1. In infinite probability theory, the entropy of a probability density function on a probability space can be defined. Given a measure space (X, \mathcal{M}, μ) consisting of a set X , a sigma algebra \mathcal{M} , and a measure μ satisfying $\mu(X) = 1$, together with a measurable function $p : X \rightarrow \mathbb{R}$ satisfying $p(x) \geq 0$ for all $x \in X$ and $\int_X p(x) d\mu(x) = 1$, the entropy is defined to be

$$(2.2) \quad H((X, \mathcal{M}, \mu, p)) := - \int_X p(x) \ln(p(x)) d\mu(x).$$

Example 2.3. In quantum mechanics of a finite number of particles, the state of a system is described by a density matrix ρ , a trace-class Hermitian operator on a separable Hilbert space \mathcal{H} . The *von Neumann entropy* of ρ is defined to be

$$(2.4) \quad H_{\text{vN}}(\rho) = -\text{tr}(\rho \ln \rho),$$

where one uses the functional calculus for operators to define $\ln \rho$. The von Neumann entropy is a measure of the entanglement of a state. However, it is not a faithful measure of the entanglement in the sense that “ $H_{\text{vN}}(\rho) > 0$ if and only if ρ is entangled” is not true [10].

Example 2.5. More generally, such as in some quantum mechanical systems of an infinite number of particles or for some quantum field theories, the state of a system is described by a normal state on a von Neumann algebra [12] (or even more generally by a state on some C^* -algebra [15]). The entropy of such a state can be defined, though the formula is a bit cumbersome to explain here and will therefore be omitted. The definition is due to Narnhofer and Thirring with earlier work by Araki [9], [1]. Note that $\mathcal{B}(\mathcal{H})$ is a von Neumann algebra and there is a one-to-one correspondence between normal states on $\mathcal{B}(\mathcal{H})$ and density matrices [7], and the formula for this entropy reduces to von Neumann’s formula.

Example 2.6. There are several notions of dynamical entropies for dynamical systems, which are usually described by some mathematical object together with a self-map. For example, a metric space (X, d) together with a uniformly continuous map $T : X \rightarrow X$ has associated with it a type of entropy whose formula is somewhat complicated and will therefore be omitted. These and closely related notions of entropy were developed by Andrey Kolmogorov, Yakov Sinai, Rufus Bowen, Efim Dinaburg, Roy Adler, Alan Konheim, and M. McAndrew [5].

Example 2.7. All of the previous examples are actually just special cases of *dynamical* entropy of C^* -algebras and von Neumann algebras. This was shown to be the case by Connes, Narnhofer, and Thirring in 1987 [4].

Perhaps the most surprising definitions of entropy, therefore, arise outside the study of dynamical systems.

Example 2.8. In Lorentzian geometry, the study of manifolds with Lorentzian metrics, some of the solutions to Einstein's equations for the metric produce black hole solutions (my apologies, but I'm actually a bit ignorant about what this means precisely so I can't give a mathematical definition here). Stephen Hawking and Demetrious Christodoulou in 1970 discovered that under any physical process, the total area of all black holes in a system never decreases, in the theory of classical Lorentzian geometry. In 1972, Jacob Bekenstein, motivated by the relationship to the second law of thermodynamics, argued that the entropy of a black hole should be proportional to its surface area [3]. However, it was not clear at all what this entropy had to do with the entropy of Gibbs or Boltzman. In 1974, Hawking studied the fluctuations of quantum fields and showed that the area of a black hole *can* decrease due to the spontaneous emission of particles, imagined as coming from an infinitesimal neighborhood of the horizon [8]. Furthermore, the entropy of this emitted radiation is proportional to the area of the black hole and the precise formula is given by

$$(2.9) \quad S_{\text{BH}} = \frac{k_B A c^3}{4G\hbar},$$

where k_B is Boltzmann's constant, G is Newton's gravitational constant, c is the speed of light, and \hbar is Planck's constant. This formula can be more concisely explained in terms of the Planck length, which is defined to be

$$(2.10) \quad \ell_{\text{P}} = \sqrt{\frac{G\hbar}{c^3}}.$$

In this case, the entropy of a black hole is given by

$$(2.11) \quad S_{\text{BH}} = \frac{k_B}{4} \left(\frac{A}{\ell_{\text{P}}^2} \right).$$

In other words, ignoring Boltzmann's constant, the entropy of a black hole is proportional to one-fourth the the area of a black hole as measured using Planck units. Just to give you an idea for how huge this number is, the entropy of a black hole of radius 1 millimeter is

$$(2.12) \quad S_{\text{BH}}(R_{\text{BH}} = 1 \text{ mm}) = \frac{k_B}{4} \left(\frac{4\pi \times 10^{-9}}{3\ell_{\text{P}}^2} \right) = 5.6 \times 10^{37}$$

Compare this to the thermodynamic entropy of a liter of water, which is

$$(2.13) \quad S_{\text{H}_2\text{O liquid}} = 3.9 \times 10^3.$$

Therefore, the amount of water needed to have the same entropy as that of a black hole of radius 1 millimeter is

$$(2.14) \quad \frac{5.6 \times 10^{37}}{3.9 \times 10^3} = 1.4 \times 10^{34}$$

liters. The radius of a spherical body of water needed to hold this many liters of water is

$$(2.15) \quad r = \sqrt[3]{\frac{3 \times 1.4 \times 10^{34}}{4\pi \times 10^3}} = 1.5 \times 10^{10}$$

meters. Compare this to the distance from the Sun to Earth which is 1.5×10^{11} meters.

In all of these examples, the entropy is often defined by some formula. Occasionally (though to the best of my knowledge, not always), there are theorems, such as Fadeev's theorem, characterizing the different types of entropies. However, all these theorems and the formulas for the entropy depend on the precise details of the mathematics surrounding them. Is there a framework that describes all such instances without referring to their technical details? Is there a general over-all perspective that says what entropy actually *is* rather than just provide a *formula* in all such instances? I do not know the answer to this question, but my personal perspective is that category theory, which disassociates itself from the details of any particular theory in mathematics, might provide some insight into these questions.

3. THE CATEGORY OF FINITE PROBABILITY SPACES

Definition 3.1. Let (X, p) and (X', p') be two finite probability spaces. A probability-preserving function from (X, p) to (X', p') is a function $f : X \rightarrow X'$ satisfying

$$(3.2) \quad p'(x') = \sum_{x \in f^{-1}(x')} p(x)$$

for all $x' \in X'$.

Example 3.3. The meaning of a probability-preserving function between probability spaces can be seen nicely in the following example. Consider the set

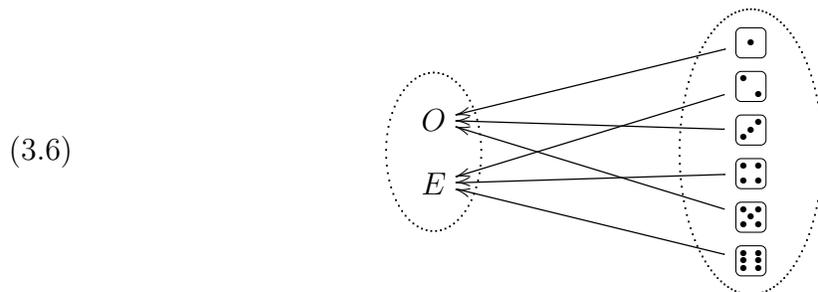
$$(3.4) \quad X := \left\{ \boxed{\bullet}, \boxed{\bullet \bullet}, \boxed{\bullet \bullet \bullet}, \boxed{\bullet \bullet \bullet \bullet}, \boxed{\bullet \bullet \bullet \bullet \bullet}, \boxed{\bullet \bullet \bullet \bullet \bullet \bullet} \right\}$$

with the probability distribution given by $\frac{1}{6}$ for each element. Consider the set

$$(3.5) \quad X' := \{O, E\}$$

consisting of just two elements (O stands for “odd” and E stands for “even”) with probability distribution given by $\frac{1}{2}$ for each element. Also consider the function $f :$

$X \longrightarrow X'$ defined by



sending \square , \square , and \square to O and sending \square , \square , and \square to E . Therefore, the probability-preserving function f is associated to the process of rolling a die and considering the likelihood of rolling an odd or an even roll versus the question of rolling a 1, 2, 3, 4, 5, or 6. There is information loss in this process because you have replaced the finer structure of the exact numbers with a coarser invariant, namely the parity of the number—odd or even. We want to quantify this information loss, and so we will seek reasonable postulates to do so.

In this regard, we must look at the *category* of finite probability spaces and probability preserving functions. The reason to emphasize category theory here is because category theory is the abstract study of processes, and if one would like to understand all kinds of entropy in terms of information change subject to some process, category theory might provide a suitable framework for such investigations.

Definition 3.7. A *category* \mathcal{C} consists of a class \mathcal{C}_0 , elements of which are referred to as *objects*, a set $\mathcal{C}_1(b, a)$, elements of which are referred to as *morphisms* from a to b , for any pair of objects a and b , a function $\mathcal{C}_1(c, b) \times \mathcal{C}_1(b, a) \longrightarrow \mathcal{C}_1(c, a)$, called the *composition*, for every triple of objects a, b , and c , and a morphism $\text{id}_a \in \mathcal{C}_1(a, a)$, called the *identity* at a , for every object a . These data are subject to several conditions.

Definition 3.8. Let **FinProb** be the category specified by the following data. **FinProb**₀ is the class of finite probability spaces (X, p) . Given two probability spaces (X, p) and (X', p') , **FinProb**₁ $((X', p'), (X, p))$ is the set of probability-preserving functions from (X, p) to (X', p') . The composition is the composition of functions. The identity at (X, p) is the identity function.

Remark 3.9. Technically, we should define **FinProb**₁ $((X', p'), (X, p))$ to be the set of a.e. equivalence classes of probability-preserving functions from (X, p) to (X', p') since one could include events with zero probability in a probability space without altering the probability distribution. This technicality is actually very important but will be ignored to a large extent for the purposes of keeping this presentation more accessible.

The Shannon entropy $H : \Delta \longrightarrow \mathbb{R}$ extends to what is known as a *functor*.

Definition 3.10. Let \mathcal{C} and \mathcal{D} be two categories. A *functor* $F : \mathcal{C} \longrightarrow \mathcal{D}$ consists of two functions $F_i : \mathcal{C}_i \longrightarrow \mathcal{D}_i$ for $i = 1, 2$ satisfying the following two conditions.

- i) For all triples of objects a, b, c in \mathcal{C} and for all pairs of composable morphisms $c \xleftarrow{g} b \xleftarrow{f} a$, the diagram

$$(3.11) \quad \begin{array}{ccc} & F_0(b) & \\ F_1(g) \swarrow & & \nwarrow F_1(f) \\ F_0(c) & \xleftarrow{F_1(g \circ f)} & F_0(a) \end{array}$$

commutes, i.e. $F_1(g \circ f) = F_1(g) \circ F_1(f)$.

- ii) For all objects a in \mathcal{C} ,

$$(3.12) \quad F(\text{id}_a) = \text{id}_{F(a)}.$$

Example 3.13. Let A be any set. From A , one can construct a category, denoted by the same letter, specified by the following data. The set of objects is given by the set A itself. The set of morphisms from a to a' is the empty set if $a \neq a'$ and is a single element set if $a = a'$. In the latter case, this element is denoted by id_a . The composition is uniquely specified. This is called the *discrete category* associated to A . A function $f : A \rightarrow B$ from a set A to a set B defines a unique functor on the associated discrete categories satisfying the condition that it agrees with f on the set of objects.

Example 3.14. The Shannon entropy function defines a functor $H : \mathbf{FinProb}_0 \rightarrow \mathbb{R}$ on the associated discrete categories by setting

$$(3.15) \quad H((X, p)) := - \sum_{x \in X} p(x) \log_2(p(x))$$

for all finite probability spaces (X, p) .

This perspective is one where information content is associated to a probability space, and not a process, which is what we have been motivated to discover. This functor satisfies many familiar conditions that will be explained in more detail later. However, it is also beneficial to instead consider the following category that is a slight modification of \mathbb{R} .

Example 3.16. Let $\mathbb{B}\mathbb{R}$ be the category specified by the following data. $\mathbb{B}\mathbb{R}_0$ is a single element set, denoted by \bullet . $\mathbb{B}\mathbb{R}_1(\bullet, \bullet)$ is the set \mathbb{R} of real numbers. The composition is addition of real numbers. The identity at \bullet is the number 0.

With this category, we can use the Shannon entropy to define a function that might potentially represent the information loss associated to probability-preserving functions.

Example 3.17. Let $F : \mathbf{FinProb} \rightarrow \mathbb{B}\mathbb{R}$ be the functor that sends a probability-preserving function $f : (X, p) \rightarrow (X', p')$ to the real number¹

$$(3.18) \quad F\left((X', p') \xleftarrow{f} (X, p)\right) := \sum_{x' \in X'} p'(x') \log_2(p'(x')) - \sum_{x \in X} p(x) \log_2(p(x)).$$

¹By definition of $\mathbb{B}\mathbb{R}$, the functor F automatically sends all probability spaces to the unique object \bullet of $\mathbb{B}\mathbb{R}$.

It is worthwhile considering special cases of this functor. First, notice that associated to every finite probability space (X, p) , there is a unique probability preserving function to $\mathbf{1}$, the probability space consisting of a single element with probability distribution equal to 1. This function sends every element of X to the unique element of $\mathbf{1}$. Denote this function by $!_{(X,p)}$. Then applying the definition of F gives

$$(3.19) \quad F\left(\mathbf{1} \xleftarrow{!_{(X,p)}} (X, p)\right) = -\sum_{x \in X} p(x) \log_2(p(x)).$$

As an additional example, notice that for the probability preserving function f of Example 3.3 that sends $\boxed{\bullet}, \boxed{\bullet\bullet}, \boxed{\bullet\bullet\bullet}$ to O and $\boxed{\bullet\bullet}, \boxed{\bullet\bullet\bullet}, \boxed{\bullet\bullet\bullet\bullet}$ to E , this gives

$$(3.20) \quad F(f) = \sum_{k=1}^2 \frac{1}{2} \log_2 \frac{1}{2} - \sum_{k=1}^6 \frac{1}{6} \log_2 \frac{1}{6} = -\log_2 2 + \log_2 6 > 0.$$

This is actually a general phenomenon, $F\left((X', p') \xleftarrow{f} (X, p)\right) \geq 0$ for all probability preserving functions between any two finite probability spaces. This follows from the following calculation

$$(3.21) \quad \begin{aligned} \sum_{x' \in X'} p'(x') \log_2 p'(x') &= \sum_{x' \in X'} \left[\left(\sum_{x \in f^{-1}(x')} p(x) \right) \log_2 \left(\sum_{\bar{x} \in f^{-1}(x')} p(\bar{x}) \right) \right] \\ &= \sum_{x' \in X'} \left[\sum_{x \in f^{-1}(x')} p(x) \log_2 \left(\sum_{\bar{x} \in f^{-1}(x')} p(\bar{x}) \right) \right] \\ &= \sum_{x' \in X'} \left[\sum_{x \in f^{-1}(x')} p(x) \log_2 \left(p(x) + \sum_{\bar{x} \in f^{-1}(x') \setminus \{x\}} p(\bar{x}) \right) \right] \\ &\geq \sum_{x' \in X'} \left[\sum_{x \in f^{-1}(x')} p(x) \log_2 p(x) \right] \\ &= \sum_{x \in X} p(x) \log_2 p(x) \end{aligned}$$

since the logarithm is a monotonically increasing function. The fact that F is a functor is quite easy to see as the value on any morphism is the difference between the Shannon entropy at its source and target.

This functor, originally described by Baez, Fritz, and Leinster [2] satisfies many other wonderful properties, but to properly describe all of them, we must understand the convex structure buried within the category **FinProb**.

4. THE CONVEX STRUCTURE ON THE CATEGORY OF FINITE PROBABILITY SPACES

Δ^{n-1} is an example of a convex set. Normally, one defines a convex set as a convex subspace of \mathbb{R}^n for some $n \in \mathbb{N}$. However, it is useful to have an abstract definition

independent of an embedding to be able to apply such a definition in a categorical context [6].

Definition 4.1. A *convex set* is a set C together with a family of functions known as *convex linear combinations* $F_\lambda : C \times C \rightarrow C$ indexed by $\lambda \in [0, 1]$ satisfying the following axioms:

$$(4.2) \quad F_0(x, y) = y \quad (\text{unit law})$$

$$(4.3) \quad F_\lambda(x, x) = x \quad (\text{idempotency})$$

$$(4.4) \quad F_\lambda(x, y) = F_{1-\lambda}(y, x) \quad (\text{parametric commutativity})$$

$$(4.5) \quad F_\lambda(F_\mu(x, y), z) = F_{\lambda \sqcup \mu}(x, F_{\lambda \sqcup \mu}(y, z)) \quad (\text{deformed parametric associativity})$$

for all $x, y, z \in C$ and $\lambda, \mu \in [0, 1]$. Here

$$(4.6) \quad \lambda \sqcup \mu := \lambda\mu \quad \& \quad \lambda \sqcup \mu := \begin{cases} \frac{\lambda(1-\mu)}{1-\lambda\mu} & \text{if } \lambda\mu \neq 1 \\ \text{arbitrary} & \text{if } \lambda = \mu = 1 \end{cases}.$$

Here “arbitrary” means that one can assign any value to the quantity.

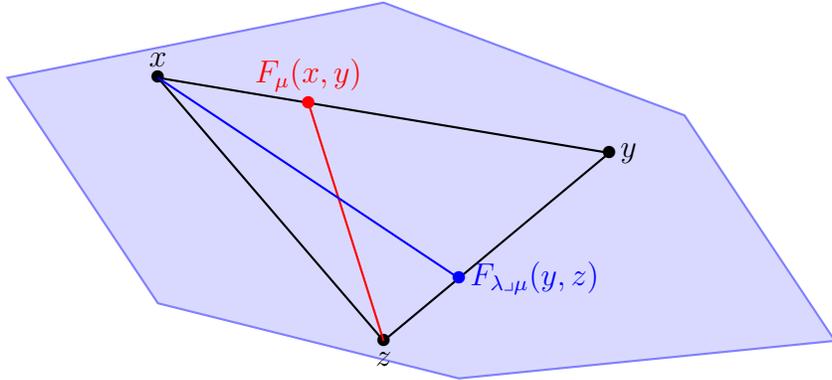
It is convenient to use the notation

$$(4.7) \quad \lambda x + (1 - \lambda)y := F_\lambda(x, y).$$

In this case, the laws take a more familiar form

$$\begin{aligned} 0x + 1y &= y \\ \lambda x + (1 - \lambda)x &= x \\ \lambda x + (1 - \lambda)y &= (1 - \lambda)y + \lambda x \\ \lambda(\mu x + (1 - \mu)y) + (1 - \lambda)z &= (\lambda \sqcup \mu)x + (1 - \lambda \sqcup \mu)((\lambda \sqcup \mu)y + (1 - \lambda \sqcup \mu)z) \end{aligned}$$

The formula for deformed parametric associativity becomes significantly more apparent when one draws a picture.



The lore of category theory is to replace equations/conditions with data, and then have those data satisfy other conditions.

Definition 4.8. Let \mathcal{C} be a category and let a, b be two objects in \mathcal{C} . An *isomorphism* from a to b , written as $b \xleftarrow{\cong} a$ is a morphism $b \xleftarrow{f} a$ for which there exists a morphism $a \xleftarrow{f^{-1}} b$ such that the diagrams

$$(4.9) \quad \begin{array}{ccc} & b & \\ f^{-1} \swarrow & & \searrow f \\ a & \xleftarrow{\text{id}_a} & a \end{array} \quad \& \quad \begin{array}{ccc} & a & \\ f \swarrow & & \searrow f^{-1} \\ b & \xleftarrow{\text{id}_b} & b \end{array}$$

both commute.

Definition 4.10. A *convex category* is a category \mathcal{C} together with a family of functors known as *convex linear combinations* $F_\lambda : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$ indexed by $\lambda \in [0, 1]$ and isomorphisms/morphisms

$$(4.11) \quad F_0(x, y) \xrightarrow{\cong} y \quad (\text{unitors})$$

$$(4.12) \quad F_\lambda(x, x) \rightarrow x \quad (\text{idempoters})$$

$$(4.13) \quad F_\lambda(x, y) \xrightarrow{\cong} F_{1-\lambda}(y, x) \quad (\text{parametric commutators})$$

$$(4.14) \quad F_\lambda(F_\mu(x, y), z) \xrightarrow{\cong} F_{\lambda \sqcup \mu}(x, F_{\lambda \sqcup \mu}(y, z)) \quad (\text{deformed parametric associators})$$

for all objects $x, y, z \in \mathcal{C}$ and for all $\lambda, \mu \in [0, 1]$. These data are subject to several conditions.

The ‘‘conditions’’ mentioned in the above theorem are referred to generically in category theory as ‘‘coherence conditions.’’ Heuristically, they turn the picture following the definition of a convex set ‘‘fuzzy’’ as if there is an additional dimension. Just as for convex sets, it is convenient to use the notation

$$(4.15) \quad \lambda x \oplus (1 - \lambda)y := F_\lambda(x, y)$$

for objects x, y in \mathcal{C} and similarly

$$(4.16) \quad \lambda f \oplus (1 - \lambda)g := F_\lambda(f, g)$$

for morphisms $x' \xleftarrow{f} x$ and $y' \xleftarrow{g} y$ in \mathcal{C} . We, however, use the direct sum notation to remind ourselves that we are summing objects and morphisms in some category and not just elements in some set. There are several examples of such structures, but we will stick to finite probability spaces.

Example 4.17. For every $\lambda \in [0, 1]$, define the convex sum F_λ on objects by

$$(4.18) \quad \begin{aligned} & \lambda(X, p) \oplus (1 - \lambda)(Y, q) := (X \amalg Y, \lambda p \oplus (1 - \lambda)q), \\ \text{where} \quad & (\lambda p \oplus (1 - \lambda)q)(z) := \begin{cases} \lambda p(z) & \text{if } z \in X \\ (1 - \lambda)q(z) & \text{if } z \in Y. \end{cases} \end{aligned}$$

The convex sum of morphisms $(X', p') \xleftarrow{f} (X, p)$ and $(Y', q') \xleftarrow{g} (Y, q)$ is given by

$$(4.19) \quad (\lambda f \oplus (1 - \lambda)g)(z) := \begin{cases} f(z) & \text{if } z \in X \\ g(z) & \text{if } z \in Y. \end{cases}$$

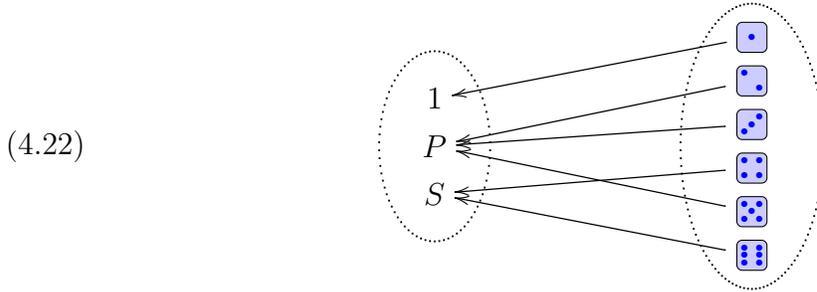
Before defining the other structure needed for a convex category, let us understand what this means in an explicit example. Let (X, p) be the probability space from Example 3.3, namely $X := \{\square, \blacksquare, \blacklozenge, \blacktriangle, \blackhexagon, \blackheptagon\}$ and $p(x) = \frac{1}{6}$ for all $x \in X$. Let $(Y, q) := (X, p)$. Fix $\lambda \in [0, 1]$. The meaning of $\lambda(X, p) \oplus (1 - \lambda)(Y, q)$ is to be interpreted as the probability space for which the events in X are distinguished from the events in Y in some way and the events in X have an additional weight of λ attached while the events in Y have an additional weight of $1 - \lambda$ attached. In other words, one first flips a coin to determine which set, X or Y , to land on and then one follows the probability distribution of the set X or Y once this action has been executed. One can therefore visualize these die as being colored

$$(4.20) \quad X \amalg Y = \{\text{red die}, \text{blue die}\}.$$

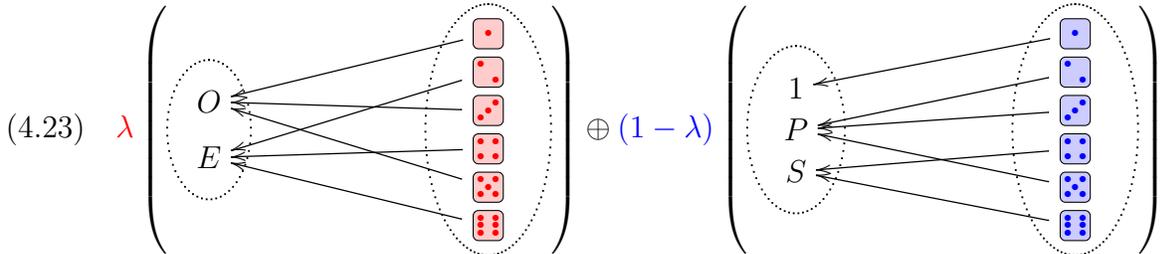
So, for example,

$$(4.21) \quad (\lambda p \oplus (1 - \lambda)q)(\text{red die}) = \frac{\lambda}{6} \quad \text{while} \quad (\lambda p \oplus (1 - \lambda)q)(\text{blue die}) = \frac{1 - \lambda}{6}.$$

Now consider the probability preserving map $f : (X, p) \rightarrow (X', p')$ from Example 3.3. Here $X' = \{O, E\}$ and p' is the uniform distribution of $\frac{1}{2}$. Consider the set $Y' := \{1, P, S\}$ (1 stands for 1, P stands for ‘‘prime,’’ and S stands for ‘‘else’’) with probability distribution $q'(1) = \frac{1}{6}$, $q'(P) = \frac{1}{2}$, and $q'(S) = \frac{1}{3}$ together with the probability preserving function $g : (Y, q) \rightarrow (Y', q')$ defined by



The function $\lambda f \oplus (1 - \lambda)g$



can be interpreted as the possibility of one of two processes to occur where the decision of which process occurs is made by flipping a λ -coin. So far, this describes the convex

linear combinations of objects and morphisms in **FinProb**. The other isomorphisms and morphisms from the definition of a convex category are left out for brevity, and we will let the reader think about what they should be (they are somewhat obvious).

Several crucial remarks are in order in regards to this example and the general structure of convex categories.

Remark 4.24.

- i) It is not true that $(X, p) \oplus (Y, q)$ is a finite probability space whenever (X, p) and (Y, q) are because the net probability of $(X, p) \oplus (Y, q)$ would be 2. This is the reason for requiring convex linear combinations.
- ii) Using probability-preserving functions as the morphisms, there is no isomorphism $0(X, p) \oplus 1(Y, q) \longrightarrow (Y, q)$ for arbitrary finite probability spaces (X, p) and (Y, q) . This is because a function cannot be an isomorphism if the cardinality of its domain does not equal the cardinality of its codomain. This is why technically the morphisms of **FinProb** should be taken to be a.e. equivalence classes of probability-preserving functions so that events with probability zero do not change the probability space.
- iii) There is in general no reasonable probability-preserving function $(X, p) \longrightarrow \lambda(X, p) \oplus (1-\lambda)(X, p)$. This is why one only asks for a morphism $\lambda(X, p) \oplus (1-\lambda)(X, p) \longrightarrow (X, p)$ as opposed to an isomorphism in the definition of a convex category. Situations for which this occurs do exist, but not for all examples of interest.

Example 4.25. Every convex set C provides a natural convex structure on the discrete category associated to C . All the morphisms/isomorphisms are identities.

Example 4.26. Let $\mathbb{B}\mathbb{R}$ be the one-object category whose morphisms are all real numbers with composition addition. Then the convex sum on objects is trivial (because there is only a single object) and the convex sum on morphisms (which are real numbers) is defined by

$$(4.27) \quad F_\lambda(a, b) := \lambda a + (1 - \lambda)b$$

for all $\lambda \in [0, 1]$ and all $a, b \in \mathbb{R}$. The other morphisms in the data of a convex category are set to be the identity (which is the number 0 in this case). In fact, $\mathbb{B}\mathbb{R}$ is a “cone category,” a concept closely related to the notion of a convex category. Essentially, this means that all objects and morphisms in the category can be scaled by any non-negative number.

5. ENTROPY IS AN AFFINE FUNCTOR

Definition 5.1. A functor $F : \mathcal{C} \longrightarrow \mathcal{D}$ is *affine* if for every $\lambda \in [0, 1]$,

$$(5.2) \quad F(\lambda x \oplus (1 - \lambda)y) = \lambda F(x) \oplus (1 - \lambda)F(y)$$

for all objects x, y in \mathcal{C} and

$$(5.3) \quad F(\lambda f \oplus (1 - \lambda)g) = \lambda F(f) \oplus (1 - \lambda)F(g)$$

for all morphisms f and g in \mathcal{C} .

Remark 5.4. Just as there are linear, convex, and concave functions on convex sets, there are even more possibilities for functors between convex categories. In particular, the subject of convex analysis is a special case of the study of convex categories and the different types of convex functors between them. This is discussed further in my thesis [11]. For the purposes of this review, we only need to focus on convex affine functors.

Example 5.5. The functor $F : \mathbf{FinProb} \rightarrow \mathbb{B}\mathbb{R}$ defined by

$$(5.6) \quad F\left((X', p') \xleftarrow{f} (X, p)\right) := \sum_{x' \in X'} p'(x') \log_2(p'(x')) - \sum_{x \in X} p(x) \log_2(p(x))$$

is affine.

Remark 5.7. Compare this result to the Shannon entropy! The Shannon entropy formula is *not* affine with respect to convex combinations of finite probability spaces! It is strictly *convex*. The deviation from the Shannon entropy being affine is precisely cancelled when one takes the difference as is done in the definition of $F : \mathbf{FinProb} \rightarrow \mathbb{B}\mathbb{R}$.

Definition 5.8. A sequence of morphisms

$$(5.9) \quad (X_n, p_n) \xrightarrow{f_n} (Y_n, q_n)$$

in $\mathbf{FinProb}$ converges to a morphism $(X, p) \xrightarrow{f} (Y, q)$ if

- (a) there exists an $N \in \mathbb{N}$ for which $X_n = X$, $Y_n = Y$, and $f_n = f$ for all $n \geq N$ and
- (b) both

$$(5.10) \quad \lim_{n \rightarrow \infty} p_n = p \quad \& \quad \lim_{n \rightarrow \infty} q_n = q,$$

where this limit is taken in the spaces,

$$(5.11) \quad \Delta^{|X|-1} \quad \& \quad \Delta^{|Y|-1},$$

respectively.

A functor $K : \mathbf{FinProb} \rightarrow \mathbb{B}\mathbb{R}$ is continuous if

$$(5.12) \quad \lim_{n \rightarrow \infty} K(f_n) = K(f)$$

whenever $\{f_n\}$ is a sequence in $\mathbf{FinProb}$ converging to f .

Example 5.13. The functor $F : \mathbf{FinProb} \rightarrow \mathbb{B}\mathbb{R}$, the difference of the Shannon entropies between the source and target, is continuous.

Theorem 5.14 (Baez-Fritz-Leinster [2]). *Let $G : \mathbf{FinProb} \rightarrow \mathbb{B}\mathbb{R}$ be a functor satisfying the following conditions.*

- i) $G(f) \geq 0$ for all morphisms f in $\mathbf{FinProb}$.
- ii) G is continuous.
- iii) $G(\mathbf{1} \xleftarrow{!2} \mathbf{2}) = 1$.
- iv) G is affine.

Then $G = F$.

The first condition associates a non-negative number to every process. G is to be interpreted as a measure of the information loss associated to a process. As we discussed above in the examples, this is what we concluded for the functor F . Functoriality reflects the fact that the information loss associated an n -step process is the sum of the information losses over all n steps. Continuity says that the information loss is perturbed in a continuous fashion when the probability spaces and processes are perturbed infinitesimally. $G(\mathbf{1} \stackrel{1}{\leftarrow} \mathbf{2}) = 1$ says that the information content of the fair coin toss is 1. Finally, the fact that G is affine says something quite intuitive. If you have two processes f and g and you toss a λ -coin to decide on which of the two processes will occur, then the information loss associated to this probabilistic process $\lambda f \oplus (1 - \lambda)g$ is the λ -convex combination of the information loss for each process.

REFERENCES

- [1] Huzihiro Araki. Relative entropy of states of von Neumann algebras. *Publ. Res. Inst. Math. Sci.*, 11(3):809–833, 1975/76.
- [2] John C. Baez, Tobias Fritz, and Tom Leinster. A characterization of entropy in terms of information loss. *Entropy*, 13(11):1945–1957, 2011.
- [3] Jacob D. Bekenstein. Black holes and entropy. *Phys. Rev. D (3)*, 7:2333–2346, 1973.
- [4] A. Connes, H. Narnhofer, and W. Thirring. Dynamical entropy of C^* algebras and von Neumann algebras. *Comm. Math. Phys.*, 112(4):691–719, 1987.
- [5] Tomasz Downarowicz. *Entropy in dynamical systems*, volume 18 of *New Mathematical Monographs*. Cambridge University Press, Cambridge, 2011.
- [6] Joe Flood. Semiconvex geometry. *J. Austral. Math. Soc. Ser. A*, 30(4):496–510, 1980/81.
- [7] Brian C. Hall. *Quantum theory for mathematicians*, volume 267 of *Graduate Texts in Mathematics*. Springer, New York, 2013.
- [8] S. W. Hawking. Particle creation by black holes. *Comm. Math. Phys.*, 43(3):199–220, 1975.
- [9] Narnhofer Heide and Walter Thirring. From relative entropy to entropy. *FZKAAA*, 17:257–265, 1985.
- [10] Elliot H. Lieb. Topics in quantum entropy and entanglement, 2014. Lectures during Princeton Summer School for Condensed Matter Physics (PSSCMP).
- [11] Arthur Parzygnat. Some 2-categorical aspects in physics, 2016. Ph.D. Thesis *CUNY Academic Works*.
- [12] Miklós Rédei and Stephen Jeffrey Summers. Quantum probability theory. *Stud. Hist. Philos. Sci. B Stud. Hist. Philos. Modern Phys.*, 38(2):390–417, 2007.
- [13] Alfréd Rényi. On measures of entropy and information. In *Proc. 4th Berkeley Sympos. Math. Statist. and Prob., Vol. I*, pages 547–561. Univ. California Press, Berkeley, Calif., 1961.
- [14] C. E. Shannon. A mathematical theory of communication. *Bell System Tech. J.*, 27:379–423, 623–656, 1948.
- [15] Robert M. Wald. *Quantum field theory in curved spacetime and black hole thermodynamics*. Chicago Lectures in Physics. University of Chicago Press, Chicago, IL, 1994.
- [16] Alfred Wehrl. General properties of entropy. *Rev. Modern Phys.*, 50(2):221–260, 1978.

MATHEMATICS DEPARTMENT, UNIVERSITY OF CONNECTICUT, STORRS, CT 06269, USA, *Email:* arthur.parzygnat@uconn.edu